



INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification⁶ :

G06F 17/30

A1

(11) International Publication Number:

WO 97/22069

(43) International Publication Date:

19 June 1997 (19.06.97)

(21) International Application Number: PCT/US96/19831

(22) International Filing Date: 10 December 1996 (10.12.96)

(30) Priority Data:

08/571,748

13 December 1995 (13.12.95)

US

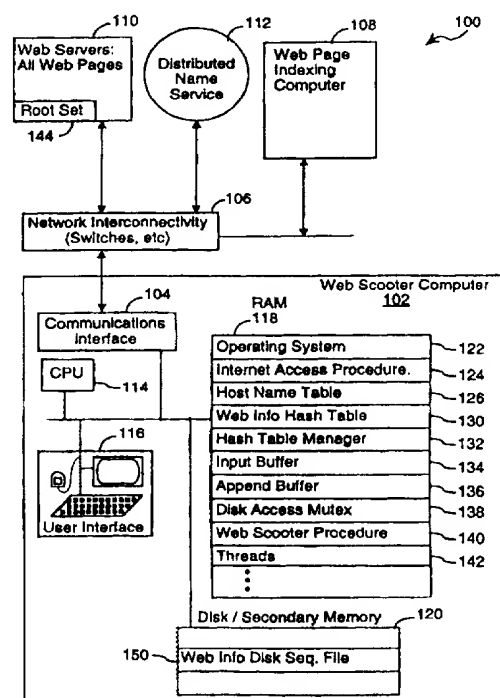
(71) Applicant: DIGITAL EQUIPMENT CORPORATION
[US/US]; 111 Powdermill Road, Maynard, MA 01754
(US).(72) Inventor: MONIER, Louis, M.; 2019 Maryland Street, Red-
wood City, CA 94061 (US).(74) Agent: NATH, Rama, B.; Digital Equipment Corporation, 111
Powdermill Road, Maynard, MA 01754 (US).(81) Designated States: AM, AT, AU, BB, BG, BR, BY, CA, CH,
CN, CZ, DE, DK, ES, FI, GB, GE, HU, JP, KE, KG, KP,
KR, KZ, LK, LT, LU, LV, MD, MG, MN, MW, NO, NZ,
PL, PT, RO, RU, SD, SE, SI, SK, TJ, TT, UA, UZ, VN,
ARIPO patent (KE, LS, MW, SD, SZ, UG), European patent
(AT, BE, CH, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU,
MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM,
GA, GN, ML, MR, NE, SN, TD, TG).

Published

*With international search report.**Before the expiration of the time limit for amending the
claims and to be republished in the event of the receipt of
amendments.*(54) Title: SYSTEM AND METHOD FOR LOCATING PAGES ON THE WORLD WIDE WEB AND FOR LOCATING DOCUMENTS
FROM A NETWORK OF COMPUTERS

(57) Abstract

A Web crawler system and method for quickly fetching and analyzing Web pages on the World Wide Web or from computers connected by a network, includes a hash table stored in a random access memory (RAM) and a sequential Web information disk file. For every Web page known to the system, the Web crawler system stores an entry in the sequential disk file as well as a smaller entry in the hash table. The hash table entry includes a fingerprint value, a fetched flag that is set true only if the corresponding Web page has been successfully fetched, and a file location indicator that indicates where the corresponding entry is stored in the sequential disk file. Each sequential disk file entry includes the URL of a corresponding Web page, plus fetch status information concerning that Web page. All accesses to the Web information disk file are made sequentially via an input buffer such that a large number of entries from the sequential disk file are moved into the input buffer as single I/O operation. The sequential disk file is then accessed from the input buffer. Similarly, all new entries to be added to the sequential file are stored in an append buffer, and the contents of the append buffer are added to the end of the sequential whenever the append buffer is filled. In this way random access to the Web information disk file is eliminated, and latency caused by disk access limitations is minimized.



FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AM	Armenia	GB	United Kingdom	MW	Malawi
AT	Austria	GE	Georgia	MX	Mexico
AU	Australia	GN	Guinea	NE	Niger
BB	Barbados	GR	Greece	NL	Netherlands
BE	Belgium	HU	Hungary	NO	Norway
BF	Burkina Faso	IE	Ireland	NZ	New Zealand
BG	Bulgaria	IT	Italy	PL	Poland
BJ	Benin	JP	Japan	PT	Portugal
BR	Brazil	KE	Kenya	RO	Romania
BY	Belarus	KG	Kyrgyzstan	RU	Russian Federation
CA	Canada	KP	Democratic People's Republic of Korea	SD	Sudan
CF	Central African Republic	KR	Republic of Korea	SE	Sweden
CG	Congo	KZ	Kazakhstan	SG	Singapore
CH	Switzerland	LI	Liechtenstein	SI	Slovenia
CI	Côte d'Ivoire	LK	Sri Lanka	SK	Slovakia
CM	Cameroon	LR	Liberia	SN	Senegal
CN	China	LT	Lithuania	SZ	Swaziland
CS	Czechoslovakia	LU	Luxembourg	TD	Chad
CZ	Czech Republic	LV	Latvia	TG	Togo
DE	Germany	MC	Monaco	TJ	Tajikistan
DK	Denmark	MD	Republic of Moldova	TT	Trinidad and Tobago
EE	Estonia	MG	Madagascar	UA	Ukraine
ES	Spain	ML	Mali	UG	Uganda
FI	Finland	MN	Mongolia	US	United States of America
FR	France	MR	Mauritania	UZ	Uzbekistan
GA	Gabon			VN	Viet Nam

**SYSTEM AND METHOD FOR LOCATING PAGES ON THE WORLD WIDE WEB
AND FOR LOCATING DOCUMENTS FROM A NETWORK OF COMPUTERS**

FIELD OF THE INVENTION

5 The present invention relates generally to systems and method for accessing documents, called pages, on the World Wide Web (WWW), and for locating documents from a network of computers, and particularly to a system and method for quickly locating and analyzing pages on the World Wide Web.

10 **BACKGROUND OF THE INVENTION**

 Web documents, herein called Web pages, are stored on numerous server computers (hereinafter "servers") that are connected to the Internet. Each page on the Web has a
15 distinct URL (universal resource locator). Many of the documents stored on Web servers are written in a standard document description language called HTML (hypertext markup language). Using HTML, a designer of Web documents can associate hypertext links or annotations with specific
20 words or phrases in a document and specify visual aspects and the content of a Web page. The hypertext links identify the URLs of other Web documents or other parts of the same document providing information related to the words or phrases.

25 A user accesses documents stored on the WWW using a Web browser (a computer program designed to display HTML documents and communicate with Web servers) running on a Web client connected to the Internet. Typically, this is done by the user selecting a hypertext link (typically
30 displayed by the Web browser as a highlighted word or phrase) within a document being viewed with the Web browser. The Web browser then issues a HTTP (hypertext transfer protocol) request for the requested document to the Web server identified by the requested document's URL.

35 In response, the designated Web server returns the requested document to the Web browser, also using the HTTP.

- 2 -

As of the end of the 1995, the number of pages on the portion of the Internet known as the World Wide Web (hereinafter the "Web") had grown several fold during the prior one year period to at least 30 million pages. The present invention is directed at a system for keeping track of pages on the Web as the Web continues to grow.

The systems for locating pages on the Web are known variously as "Web crawlers," "Web spiders" and "Web robots." The present invention has been coined a "Web scooter" because it is so much faster than all known Web crawlers. The terms "Web crawler," "Web spider," "Web scooter," "Web crawler computer system," and "Web scooter computer system" are used interchangeably in this document.

Prior art Web crawlers work generally as follows. Starting with a root set of known Web pages, a disk file is created with a distinct entry for every known Web page. As additional Web pages are fetched and their links to other pages are analyzed, additional entries are made in the disk file to reference Web pages not previously known to the Web crawler. Each entry indicates whether or not the corresponding Web page has been processed as well as other status information. A Web crawler processes a Web page by (A) identifying all links to other Web pages in the page being processed and storing related information so that all of the identified Web pages that have not yet been processed are added to a list of Web pages to be processed or other equivalent data structure, and (B) passing the Web page to an indexer or other document processing system.

The information about the Web pages already processed is generally stored in a disk file, because the amount of information in the disk file is too large to be stored in random access memory (RAM). For example, if an average of 100 bytes of information are stored for each Web page entry, a data file representing 30 million Web pages would

occupy about 3 Gigabytes, which is too large for practical storage in RAM.

Next we consider the disk I/O incurred when processing one Web page. For purposes of this discussion we will
5 assume that a typical Web page contains 20 references to other Web pages, and that a disk storage device can handle no more than 50 seeks per second. The Web crawler must evaluate each of the 20 page references in the page being processed to determine if it already knows about those
10 pages. To do this it must attempt to retrieve 20 records from the Web information disk file. If the record for a specified page reference already exists, then that reference is discarded because no further processing is needed. However, if a record for a specified page is not
15 found, an attempt must be made to locate a record for each possible alias of the page's address, thereby increasing the average of number of disk record seeks needed to analyze an average Web page to about 50 disk seeks per page.

20 If a disk file record for a specified page reference does not already exist a new record for the referenced page is created and added to the disk file, and that page reference is either added to a queue of pages to be processed, or the disk file entry is itself used to
25 indicate that the page has not yet been fetched and processed.

Thus, processing a single Web page requires approximately 20 disk seeks (for reading existing records and for writing new records). As a result, given a
30 limitation of 50 disk seeks per second, only about one Web page can be processed per second.

In addition, there is a matter of network access latency. On average, it takes about 3 seconds to retrieve a Web page, although the amount of time is highly variable
35 depending on the location of the Web server and the particular hardware and software being used on both the Web

server and on the Web crawler computer. Network latency thus also tends to limit the number Web pages that can be processed by prior art Web crawlers to about 0.33 Web pages per second. Due to disk "seek" limitations, network
5 latency, and other delay factors, a typical prior art Web crawler cannot process more than about 30,000 Web pages per day.

Due to the rate at which Web pages are being added to the Web, and the rate at which Web pages are being deleted
10 and revised, processing 30,000 Web pages per day is inadequate for maintaining a truly current directory or index of all the Web pages on the Web. Ideally, a Web crawler should be able to visit (i.e., fetch and analyze) at least 2.5 million Web pages per day.

It is therefore desirable to have a Web crawler with
15 such high speed capacity. An object of the present invention to provide an improved Web crawler that can process millions of Web pages per day. It is a related goal of the present invention to provide an improved Web
20 crawler that overcomes the aforementioned disk "seek" limitations and network latency limitations so as to enable the Web crawler's speed of operation to be limited primarily only by the processing speed of the Web crawler's CPU. It is yet another related goal of the present
25 invention to provide a Web crawler system than can fetch and analyze, on average, at least 30 Web pages per second, and more preferably at least 100 Web pages per second.

SUMMARY OF THE INVENTION

30 The invention, in its broad form, resides in a system for locating Web pages as recited in claim 1 and a method for locating Web pages as recited in claim 6.

Described hereinafter is a system and method for
quickly locating and making a directory of Web pages on the
35 World Wide Web. The Web crawler system includes a hash table stored in random access memory (RAM) and a sequential

- 5 -

file (herein called the "sequential disk file" or the "Web information disk file") stored in secondary memory, typically disk storage. For every Web page known to the system, the Web crawler system stores an entry in the sequential disk file as well as a smaller entry in the hash table. The hash table entry includes a fingerprint value, a fetched flag that is set true only if the corresponding Web page has been successfully fetched, and a file location indicator that indicates where the corresponding entry is stored in the sequential disk file. Each sequential disk file entry includes the URL of a corresponding Web page, plus fetch status information concerning that Web page.

All accesses to the Web information disk file are made sequentially via an input buffer such that a large number of entries from the sequential disk file are moved into the input buffer as single I/O operation. The sequential disk file is then accessed from the input buffer. Similarly, all new entries to be added to the sequential file are stored in an append buffer, and the contents of the append buffer are added to the end of the sequential disk file whenever the append buffer is filled. In this way random access to the Web information disk file is eliminated, and latency caused by disk access limitations is minimized.

The procedure for locating and processing Web pages includes sequentially reviewing all entries in the sequential file and selecting a next entry that meets with established selection criteria. When selecting the next file entry to process, the hash table is checked for all known aliases of the current entry candidate to determine if the Web page has already been fetched under an alias. If the Web page has been fetched under an alias, the error type field of the sequential file entry is marked as a "non-selected alias" and the candidate entry is not selected.

Once a next Web page reference entry has been selected, the Web crawler system attempts to fetch the

- 6 -

corresponding Web page. If the fetch is unsuccessful, the fetch status information in the sequential file entry for that Web page is marked as a fetch failure in accordance with the error return code returned to the Web crawler. If the fetch is successful, the fetch flag in the hash table entry for the Web page is set, as is a similar fetch flag in the sequential disk file entry (in the input buffer) for the Web page. In addition, each URL link in the fetched Web page is analyzed. If an entry for the URL referenced by the link or for any defined alias of the URL is already in the hash table, no further processing of the URL link is required. If no such entry is found in the hash table, the URL represents a "new" Web page not previously included in the Web crawler's database of Web pages and therefore an entry for the new Web page is added to the sequential disk file (i.e., it is added to the portion of the disk file in the append buffer). The new disk file entry includes the URL referenced by the link being processed, and is marked "not fetched". In addition, a corresponding new entry is added to the hash table, and the fetch flag of that entry is cleared to indicate that the corresponding Web page has not yet been fetched. In addition to processing all the URL links in the fetched page, the Web crawler sends the fetched page to an indexer for further processing.

BRIEF DESCRIPTION OF THE DRAWINGS

A more detailed understanding of the invention may be had from the following description of a preferred embodiment, given by way of example, and to be understood in conjunction with the accompanying drawings, wherein:

- ♦ Fig. 1 is a block diagram of a preferred embodiment of a Web crawler system in accordance with a preferred embodiment of the present invention.

- 7 -

♦ Fig. 2 is a block diagram of the hash table mechanism used in a preferred embodiment of the present invention.

5 ♦ Fig. 3 is a block diagram of the sequential Web information disk file and associated data structures used in a preferred embodiment of the present invention.

10 ♦ Fig. 4 is a flow chart of the Web crawler procedure used in a preferred embodiment of the present invention.

DESCRIPTION OF THE PREFERRED EMBODIMENTS

Referring to Fig. 1, there is shown a distributed computer system 100 having a Web scooter computer system 102. The Web scooter is connected by a communications interface 104 and a set of Internet and other network connections 106 to the Internet and a Web page indexing computer 108. In some embodiments the Web page indexing computer 108 is coupled directly to the Web scooter 102 through a private communication channel, without the use of a local or wide area network connection. The portions of the Internet to which the Web scooter 102 is connected are (A) Web servers 110 that store Web pages, and (B) servers that cooperate in a service known as the Distributed Name Service (DNS) collectively referenced here by reference numeral 112. For the purposes of this document it can be assumed that the DNS 112 provides any requester with the set of all defined aliases for any Internet host name, and that Internet host names and their aliases form a prefix portion of every URL.

30 In the preferred embodiment, the Web scooter 102 is an Alpha workstation computer made by Digital Equipment Corporation; however virtually any type of computer can be used as the Web scooter computer. In the preferred embodiment, the Web scooter 102 includes a CPU 114, the

previously mentioned communications interface 104, a user interface 116, random access memory (RAM) 118 and disk memory (disk) 120. In the preferred embodiment the communications interface 104 is a very high capacity communications interface that can handle 1000 or more overlapping communication requests with an average fetch throughput of at least 30 Web pages per second.

In the preferred embodiment, the Web scooter's RAM has a Gigabyte of random access memory and stores:

- a multitasking operating system 122;
- an Internet communications manager program 124 for fetching Web pages as well as for fetching alias information from the DNS 112;
- a host name table 126, which stores information representing defined aliases for host names;
- a Web information hash table 130;
- a hash table manager procedure 132;
- an input buffer 134 and an append buffer 136;
- a mutex 138 for controlling access to the hash table 130, input buffer 134 and append buffer 136;
- a Web scooter procedure 140; and
- thread data structures 142 for defining T1 threads of execution, where the value of T1 is an integer selectable by the operator of the Web scooter computer system 102 (e.g., T1 is set at a value of 1000 in the preferred embodiment).

Disk storage 120 stores a Web information disk file 150 that is sequentially accessed through the input buffer 134 and append buffer 136, as described in more detail below.

The host name table 126 stores information representing, among other things, all the aliases of each host name that are known to the DNS 112. The aliases are effectively a set of URL prefixes which are substituted by the Web scooter procedure 140 for the host name portion of

a specified Web page's URL to form a set of alias URLs for the specified Web page.

The use and operation of the above mentioned data structures and procedures will next be described with reference to Figs. 1 through 4 and with reference to Tables 1 and 2. Tables 1 and 2 together contain a pseudocode representation of the Web scooter procedure. While the pseudocode employed here has been invented solely for the purposes of this description, it utilizes universal computer language conventions and is designed to be easily understandable by any computer programmer skilled in the art.

Web Information Hash Table

Referring to Fig. 2, the Web information hash table includes a distinct entry 160 for each Web page that has been fetched and analyzed by the Web scooter system as well as each Web page referenced by a URL link in a Web page that has been fetched and analyzed. Each such entry includes:

- a fingerprint value 162 that is unique to the corresponding Web page;
- a one bit "fetched flag" 164 that indicates whether or not the corresponding Web page has been fetched and analyzed by the Web scooter; and
- a file location value 166 that indicates the location of a corresponding entry in the Web information disk file 150.

In the preferred embodiment, each fingerprint value is 63-bits long, and the file location values are each 32-bits long. As a result each hash table entry 160 occupies exactly 12 bytes in the preferred embodiment. While the exact size of the hash table entries is not important, it is important that each hash table entry 160 is

- 10 -

significantly smaller (e.g., at least 75% smaller on average) than the corresponding disk file entry.

The hash table manager 132 receives, via its "interface" 170, two types of procedure calls from the Web scooter procedure 140:

- a first request asks the hash table manager 132 whether or not an entry exists for a specified URL, and if so, whether or not the fetched flag of that record indicates that the corresponding Web page has previously been fetched and analyzed; and
- a second request asks the hash table manager to store a new entry in the hash table 130 for a specified URL and a specified disk file location.

The hash table manager 132 utilizes a fingerprint hash function 172 to compute a 63-bit fingerprint for every URL presented to it. The fingerprint function 172 is designed to ensure that every unique URL is mapped into a similarly unique fingerprint value. The fingerprint function generates a compressed encoding of any specified Web page's URL. The design of appropriate fingerprint functions is understood by persons of ordinary skill in the art. It is noted that while there are about 2^{25} to 2^{26} Web pages, the fingerprints can have 2^{63} distinct values.

When the Web scooter procedure 140 asks the hash table manager 132 whether or not the hash table already has an entry for a specified URL, the hash table manager (A) generates a fingerprint of the specified URL using the aforementioned fingerprint hash function 172, (B) passes that value to a hash table position function 174 that determines where in the hash table 130 an entry having that fingerprint value would be stored, (C) determines if such an entry is in fact stored in the hash table, (D) returns a failure value (e.g., -1) if a matching entry is not found, and (E) returns a success value (e.g., 0) and fetched flag

- 11 -

value and disk position value of the entry if the entry is found in the hash table.

In the preferred embodiment, the hash table position function 174 determines the position of a hash table entry based on a predefined number of low order bits of the fingerprint, and then follows a chain of blocks of entries for all fingerprints with the same low order bits. Entries 160 in the hash table 130 for a given value of the low order bits are allocated in blocks of B1 entries per block, where B1 is a tunable parameter. The above described scheme used in the preferred embodiment has the advantage of storing data in a highly dense manner in the hash table 130. As will be understood by those skilled in the art, many other hash table position functions could be used.

When the Web scooter procedure 140 asks the hash table manager 132 to store a new hash table entry for a specified URL and a specified disk file location, the hash table manager (A) generates a fingerprint of the specified URL using the aforementioned fingerprint hash function 172, (B) passes that value to a hash table position function 174 that determines where in the hash table 130 an entry having that fingerprint value should be stored, and (C) stores a new entry 160 in the hash table at the determined position, with a fetch flag value that indicates the corresponding Web page has not yet been fetched, and also containing the fingerprint value and the specified disk file position.

Web Information Disk File and Buffers

Referring to Fig. 3 and Table 2, disk access operations are minimized through the use of an input buffer 134 and an append buffer 136, both of which are located in RAM. Management of the input and append buffers is performed by a background sequential disk file and buffer handler procedure, also known as the disk file manager.

In the preferred embodiment, the input buffer and append buffer are each 50 to 100 Megabytes in size. The input buffer 134 is used to store a sequentially ordered contiguous portion of the Web information disk file 150. The Web scooter procedure maintains a pointer 176 to the next entry in the input buffer to be processed, a pointer 178 to the next entry 180 in the Web information disk file 150 to be transferred to the input buffer 134, as well as a number of other bookkeeping pointers required for coordinating the use of the input buffer 134, append buffer 136 and disk file 150.

All accesses to the Web information disk file 150 are made sequentially via the input buffer 134 such that a large number of entries from the sequential disk file are moved into the input buffer as single I/O operation. The sequential disk file 150 is then accessed from the input buffer. Similarly, all new entries to be added to the sequential file are stored in the append buffer 136, and the contents of the append buffer are added to the end of the sequential whenever the append buffer is filled. In this way random access to the Web information disk file is eliminated, and latency caused by disk access limitations is minimized.

Each time all the entries in the input buffer 134 have been scanned by the Web scooter, all updates to the entries in the input buffer are stored back into the Web information disk file 150 and all entries in the append buffer 136 are appended to the end of the disk file 150. In addition, the append buffer 136 is cleared and the next

- 13 -

set of entries in the disk file, starting immediately after the last set of entries to be copied into the input buffer 134 (as indicated by pointer 178), are copied into the input buffer 134. When the last of the entries in the disk file have been scanned by the Web scooter procedure, scanning resumes at the beginning of the disk file 150.

Whenever the append buffer 136 is filled with new entries, its contents are appended to the end of the disk file 150 and then the append buffer is cleared to receive new entries.

Each entry 180 in the Web information disk file 150 stores:

- a variable length URL field 182 that stores the URL for the Web page referenced by the entry;
- a fetched flag 184 that indicates whether or not the corresponding Web page has been fetched and analyzed by the Web scooter;
- a timestamp 186 indicating the date and time the referenced Web page was fetched, analyzed and indexed;
- a size value 188 indicating the size of the Web page;
- an error type value 190 that indicates the type of error encountered, if any, the last time an attempt was made to fetch the referenced Web page or if the entry represents a duplicate (i.e., alias URL) entry that should be ignored; and
- other fetch status parameters 192 not relevant here.

Because the URL field 182 is variable in length, the records 180 in the Web information disk file 150 are also variable in length.

Web Scooter Procedure

Referring now to Figs. 1-4 and the pseudocode in Table 1, the Web scooter procedure 140 in the preferred embodiment works as follows. When the Web scooter

procedure begins execution, it initializes (200) the system's data structures by:

- 5 • scanning through a pre-existing Web information disk file 150 and initializing the hash table 130 with entries for all entries in the sequential disk file;
- copying a first batch of sequential disk entries from the disk file 150 into the input buffer 134;
- defining an empty append buffer 136 for new sequential
- 10 file entries; and
- defining a mutex 138 for controlling access to the input buffer 134, append buffer 136 and hash table 130.

15 The Web scooter initializer then launches T1 threads (e.g., 1000 threads are launched in the preferred embodiment), each of which executes the same scooter procedure.

20 The set of entries in the pre-existing Web information disk file 150, prior to execution of the Web scooter initializer procedure, is called the "root set" 144 of known Web pages. The set of "accessible" Web pages consists of all Web pages referenced by URL links in the root set and all Web pages referenced by URL links in other accessible Web pages. Thus it is possible that some Web

25 pages are not accessible to the Web scooter 102 because there are no URL link connections between the root set and those "inaccessible" Web pages.

30 When information about such Web pages becomes available via various channels, the Web information disk file 150 can be expanded (thereby expanding the root set 144) by "manual" insertion of additional entries or other mechanisms to include additional entries so as to make accessible the previously inaccessible Web pages.

35 The following is a description of the Web scooter procedure executed by all the simultaneously running threads. The first step of the procedure is to request and

- 15 -

wait for the mutex (202). Ownership of the mutex is required so that no two threads will process the same disk file entry, and so that no two threads attempt to write information at the same time to the hash table, input
5 buffer, append buffer or disk file. The hash table 130, input buffer 134, append buffer 136 and disk file 150 are herein collectively called the "protected data structures," because they are collectively protected by use of the mutex. Once a thread owns the mutex, it scans the disk
10 file entries in the input buffer, beginning at the next entry that has not yet been scanned (as indicated by pointer 176), until it locates and selects an entry that meets defined selection criteria (204).

For example, the default selection criteria is: any
15 entry that references a Web page denoted by the entry as never having been fetched, or which was last fetched and analyzed more than H1 hours ago, where H1 is a operator selectable value, but excluding entries whose error type field indicates the entry is a duplicate entry (i.e., a
20 "non-selected alias," as explained below). If H1 is set to 168, all entries referencing Web pages last fetched analyzed more than a week ago meet the selection criteria. Another example of a selection criteria, in which Web page size is taken into account, is: an entry representing a Web
25 page that has never been fetched, or a Web page of size greater than S1 that was last fetched and analyzed more than H1 hours ago, or a Web page of size S1 or less that was last fetched and analyzed more than H2 hours ago, but excluding entries whose error type field indicates the
30 entry is a "non-selected alias," where S1, H1 and H2 are operator selectable values.

When selecting the next entry to process, the hash table is checked for all known aliases of the current entry candidate to determine if the Web page has already been
35 fetched under an alias. In particular, if an entry meets the defined selection criteria, all known aliases of the

- 16 -

URL for the entry are generated using the information in the host name table 126, and then the hash table 130 is checked to see if it stores an entry for any of the alias URLs with a fetched flag that indicates the referenced Web page has been fetched under that alias URL. If the Web page referenced by the current entry candidate in the input buffer is determined to have already been fetched under an alias URL, the error type field 190 of that input buffer entry is modified to indicate that this entry is a "non-selected alias," which prevents the entry from being selected for further processing both at this time and in the future.

Once a Web page reference entry has been selected, the mutex is released so that other threads can access the protected data structures (206). Then the Web scooter procedure attempts to fetch the corresponding Web page (208). After the fetch completes or fails the procedure once again requests and waits for the mutex (210) so that it can once again utilize the protected data structures.

If the fetch is unsuccessful (212-N), the fetch status information in the sequential file entry for that Web page is marked as a fetch failure in accordance with the error return code returned to the Web crawler (214). If the fetch is successful (212-Y), the fetch flag 164 in the hash table entry 160 for the Web page is set, as is the fetch flag 184 in the sequential disk file entry 180 (in the input buffer) for the Web page. In addition, each URL link in the fetched Web page is analyzed (216).

After the fetched Web page has been analyzed, or the fetch failure has been noted in the input buffer entry, the mutex is released so that other threads can access the protected data structures (218).

The procedure for analyzing the URL links in the fetched Web page is described next with reference to Fig. 4B. It is noted here that a Web page can include URL links to documents, such as image files, that do not contain

- 17 -

information suitable for indexing by the indexing system 108. These referenced documents are often used as components of the Web page that references them. For the purposes of this document, the URL links to component files such as image files and other non-indexable files are not "URL links to other Web pages." These URL links to non-indexable files are ignored by the Web scooter procedure.

Once all the URL links to other Web pages have been processed (230), the fetched Web page is sent to the indexer for indexing (232) and the processing of the fetched Web page by the Web scooter is completed. Otherwise, a next URL link to a Web page is selected (234). If there is already a hash table entry for the URL associated with the selected link (236), no further processing of that link is required and a next URL link is selected (234) if there remain any unprocessed URL links in the Web page being analyzed.

If there isn't already a hash table entry for the URL associated with the selected link (236), all known aliases of the URL for the entry are generated using the information in the host name table 126, and then the hash table 130 is checked to see if it stores an entry for any of the alias URLs (238). If there is an entry in the hash table for any of the alias URLs, no further processing of that link is required and a next URL link is selected (234) if there remain any unprocessed URL links in the Web page being analyzed.

If no entry is found in the hash table for the selected link's URL or any of its aliases, the URL represents a "new" Web page not previously included in the Web crawler's database of Web pages and therefore an entry for the new Web page is added to the portion of the disk file in the append buffer (240). The new disk file entry includes the URL referenced by the link being processed, and is marked "not fetched". In addition, a corresponding

new entry is added to the hash table, and the fetch flag of that entry is cleared to indicate that the corresponding Web page has not yet been fetched (240). Then processing of the Web page continues with the next unprocessed URL link in the Web page (234), if there remain any unprocessed URL links in the Web page.

The Web information hash table 130 is used, by procedures whose purpose and operation are outside the scope of this document, as an index into the Web information disk file 150 because the hash table 130 includes disk file location values for each known Web page.

In other words, an entry in the Web information disk file is accessed by first reading the disk file address in the corresponding entry in the Web information hash table and then reading the Web information disk file entry at that address.

Alternative Embodiments

Any data structure that has the same properties of the Web information hash table 130, such as a balanced tree, a skip list, or the like, could be used in place of the hash table structure 130 of the preferred embodiment.

As a solution, the present invention uses three primary mechanisms to overcome the speed limitations of prior art Web crawlers.

First, a Web page directory table is stored in RAM with sufficient information to determine which Web pages links represent new Web pages not previously known to the Web crawler, enabling received Web pages to be analyzed without having to access a disk file.

Second, a more complete Web page directory is accessed only in sequential order, and performing those accesses via large input and append buffers that reduce the number of disk accesses performed to the point that disk accesses do not have a significant impact on the speed performance of the Web crawler.

Third, by using a large number of simultaneously active threads to execute the Web scooter procedure, and by providing a communications interface capable of handling a similar number of simultaneous communication channels to Web servers, the present invention avoids the delays caused by network access latency.

In particular, while numerous ones of the threads are waiting for responses to Web page fetch requests, other ones of the threads are analyzing received Web pages. By using a large number of threads all performing the same Web scooter procedure, there will tend to be, on average, a queue of threads with received Web pages that are waiting for the mutex so that they can process the received Web pages. Also, the Web page fetches will tend to be staggered over time. As a result, the Web scooter is rarely in a state where it is waiting to receive a Web page and has no other work to do. Throughput of the Web scooter can then be further increased by using a multiprocessor workstation and further increasing the number of threads that are simultaneously executing the Web scooter procedure.

While the present invention has been described with reference to a few specific embodiments, the description is illustrative of the invention and is not to be construed as limiting the invention. Various modifications may be made without departing from the scope of the invention as presented and claimed herein.

TABLE 1**Pseudocode Representation of Web Scooter Procedure**

Procedure: Web Scooter

```
5  {
    /* Initialization Steps */
    Scan through pre-existing Web information disk file and
        initialize Hash Table with entries for all entries in
        the sequential file
10  Read first batch of sequential disk entries into input
        buffer in RAM
    Define empty Append Buffer for new sequential file entries
    Define Mutex for controlling access to Input Buffer, Append
        Buffer and Hash Table
15  Launch 1000 Threads, each executing same Scooter Procedure
    }
```

Procedure: Scooter

```
    {
20  Do Forever:
        {
            Request and Wait for Mutex
            Read sequential file (in Input Buffer) until a new URL
                to process is selected in accordance with established
25  URL selection criteria. When selecting next URL to
                process, check Hash Table for all known aliases of URL
                to determine if the Web page has already been fetched
                under an alias, and if the Web page has been fetched
                under an alias mark the Error Type field of the
30  sequential file entry as a "non-selected alias."
            /* Example of Selection Criteria: URL has never been
                fetched, or was last fetched more than H1
                hours ago, and is not a non-selected alias */
            Release Mutex
```

- 21 -

```
Fetch selected Web page
Request and Wait for Mutex
If fetch is successful
{
5   Mark page as fetched in Hash Table and Sequential File
    entry in Input Buffer
    /* Analyze Fetched Page */
    For each URL link in the page
    {
10   If URL or any defined alias is already in the Hash
        Table
        {      Do Nothing      }
    Else
        {
15   /* the URL represents a "New" Web Page not
            previously included in the database */
            Add new entry for corresponding Web page to the
            Append Buffer, with entry marked "not fetched"
            Add entry to Hash Table, with entry marked "not
20   fetched"
        }
    }
    Send Fetched Page to Indexer for processing
}
25 Else
    {
        Mark the entry in Input Buffer currently being
        processed with appropriate "fetch failure" error
        indicator based on return code received
30   }
    Release Mutex
} /* End of Do Forever Loop */
}
```

- 22 -

TABLE 2

Pseudocode Representation for Background
Sequential File Buffer Handler

```
5  Procedure: Background Sequential File Buffer Handler (a/k/a
    the disk file manager)
    {
    Whenever a "read sequential file" instruction overflows the
        Input Buffer
10   {
        Copy the Input Buffer back to the sequential disk file
        Read next set of entries into Input Buffer
        Append contents of Append Buffer to the end of the
            sequential disk file
15   Clear Append Buffer to prepare for new entries
    }

    Whenever an "add entry to sequential file" causes the
        Append Buffer to Overflow
20   {
        Append contents of Append Buffer to the end of the
            sequential disk file
        Clear Append Buffer to prepare for new entries
        Add pending new entry to the beginning of the Append
25   Buffer
    }
}
```

WHAT IS CLAIMED IS:

1 1. A system for locating data sets comprising Web pages
2 stored on remotely located accessible computers, each Web
3 page having a unique URL (universal resource locator), at
4 least some of said Web pages including URL links to other
5 ones of the Web pages, the system comprising:

6 a communications interface for fetching specified
7 ones of the Web pages from said remotely located
8 computers in accordance with corresponding URLs;

9 a Web information file having a set of entries,
10 each entry denoting, for a corresponding Web page, a
11 URL and fetch status information;

12 a Web information table, stored in RAM (random
13 access memory), having a set of entries, each entry
14 denoting a fingerprint value and fetch status
15 information for a corresponding Web page; and

16 means to implement a Web scooter procedure,
17 executed by the system, for fetching and analyzing Web
18 pages, said Web scooter procedure including
19 instructions for fetching Web pages whose Web
20 information file entries meet predefined selection
21 criteria based on said fetch status information, for
22 determining for each URL link in each received Web
23 page whether a corresponding entry already exists in
24 the Web information table, and for each URL link which
25 does not have a corresponding entry in the Web
26 information table adding a new entry in the Web
27 information table and a corresponding new entry in the
28 Web information file.

1 2. The system of claim 1, including multiple threads that
2 each execute the Web scooter procedure during overlapping
3 time periods, including means such that while some of the
4 threads are fetching Web pages, other ones of the Web pages
5 are analyzing fetched Web pages.

1 3. The system of claim 2, including a mutex, wherein said
2 Web scooter procedure executed by each of the threads
3 includes instructions for requesting and waiting for the
4 mutex before accessing the Web information table and Web
5 information file.

1 4. The system of claim 3, including:

2 an input buffer and an append buffer;

3 a file manager for storing blocks of sequentially
4 ordered entries from the Web information file into the
5 input buffer;

6 said Web scooter procedure scanning and analyzing
7 Web information file entries in the input buffer to
8 locate said Web information file entries that meet
9 said predefined selection criteria;

10 said Web scooter procedure storing in said append
11 buffer all entries to be added to Web information said
12 file; and

13 said file manager for moving multiple entries in
14 the append buffer to the Web information file.

1 5. The system of claim 1, wherein each of the entries in
2 the second memory include an address of a corresponding
3 entry in the first memory.

1 6. A method of locating data sets comprising Web pages
2 stored on remotely located but accessible computers, each
3 Web page having a unique URL (universal resource locator),
4 at least some of said Web pages including URL links to
5 other ones of the Web pages, comprising the steps of:

6 storing a Web information file having a set of
7 entries, each entry denoting, for a corresponding Web
8 page, a URL and fetch status information;

9 storing in RAM (random access memory) a Web
10 information table having a set of entries, each entry
11 denoting a fingerprint value and fetch status
12 information for a corresponding Web page; and

13 executing a Web scooter procedure, for fetching
14 and analyzing Web pages, including (A) sequentially
15 scanning entries in the Web information file do
16 determine which of said entries meet predefined
17 selection criteria, (B) fetching Web pages whose Web
18 information file entries meet said predefined
19 selection criteria, (C) determining for each URL link
20 to another Web page in each received Web page whether
21 a corresponding entry already exists in the Web
22 information table, and (D) for each URL link which
23 does not have a corresponding entry in the Web
24 information table adding a new entry in the Web
25 information table and a corresponding new entry in the
26 Web information file.

1 7. The method of claim 6, including executing said Web
2 scooter procedure in multiple threads during overlapping
3 time periods, such that while some of the threads are
4 fetching Web pages, other ones of the Web pages are
5 analyzing fetched Web pages.

1 8. The method of claim 7, including:
2 defining a mutex;
3 while executing said Web scooter procedure in
4 each of said threads, requesting and waiting for the
5 mutex before accessing the Web information table and
6 Web information file.

1 9. The method of claim 8, including the steps of:
2 defining an "input buffer" and an "append buffer"
3 in said RAM;
4 storing blocks of sequentially ordered entries
5 from the Web information file into the input buffer;
6 said step of sequentially scanning entries in the
7 Web information file comprising the step of including
8 scanning the Web information file entries in the input

- 26 -

9 buffer to determine which of said Web information file
10 entries meet said predefined selection criteria;
11 storing in said append buffer all entries to
12 be added to said file; and
13 moving multiple entries in the append buffer to
14 the Web information file.

1 10. The method of claim 6, wherein each of the entries in
2 the Web information table includes an address of a
3 corresponding entry in the Web information file, said
4 method including:

5 accessing one of said entries in said Web
6 information file by reading the address in a
7 corresponding ones of the entries in the Web
8 information table and then reading said one entry in
9 said Web information file at said address.

1 11. An apparatus for locating data sets stored on
2 computers connected by a network, each data set being
3 uniquely identified by an address, at least some of the
4 data sets including one or more linked addresses of other
5 data sets stored on the computers, comprising:

6 a communications interface connected to the
7 network for sending requests to the computers for
8 identified ones of the data sets, each request
9 including the address of the identified one of the
10 data sets, and for receiving data sets in response to
11 said requests;

12 a first memory storing a first set of entries,
13 each entry of the first set including the address of a
14 corresponding data set and status information for the
15 corresponding data set;

16 a second memory storing a second set of entries,
17 each entry of the second set including an encoding of
18 the address of a corresponding data set and an

19 encoding of status information for the corresponding
20 data set; and

21 thread means, coupled to the first and second
22 memories and to the communications interface, for
23 sequentially reading the entries of the first set,
24 generating the requests for those identified ones of
25 the data sets that have corresponding entries in the
26 first set that meet predefined status-based selection
27 criteria, and, in response to receiving the identified
28 data sets, creating new entries in said first and
29 second sets corresponding to each of at least a subset
30 of the addresses in the received data sets for which
31 there is no corresponding entry in the second set.

1 12. The apparatus of claim 11, wherein each of the entries
2 in the second set include an address of a corresponding
3 entry in the first set, said second set of entries being
4 for indexing the first set of entries.

1 13. The apparatus of claim 11, including a multiplicity of
2 said thread means such that while some of the thread means
3 are generating said requests and receiving said identified
4 data sets, other ones of the thread means are creating new
5 entries in said first and second memories.

1 14. The apparatus of claim 13, including a mutex, wherein
2 each of said thread means includes logic for requesting and
3 waiting for the mutex before accessing the first memory and
4 second memory.

1 15. The apparatus of claim 14, includes:
2 an input buffer and an append buffer, located in
3 said second memory;
4 a manager that stores in the input buffer
5 sequentially ordered groups of the entries in the
6 first memory;

- 28 -

each of said thread means including means for scanning and analyzing entries in the input buffer to locate said entries that meet said predefined status-based selection criteria; and

each of said thread means storing in said append buffer all entries to be added to said first memory;

said manager also having means for moving multiple entries in the append buffer to the first memory.

16. A method of locating data sets stored on computers connected by a network, each data set being uniquely identified by an address, at least some of said data sets including one or more linked addresses of other data sets stored on the computers, comprising the steps of:

(A) storing in a first memory a first set of entries, each entry in said first set including the address of a corresponding data set and status information for the corresponding data set;

(B) storing in a second memory a second set of entries, each entry in said second set including an encoding of the address of a corresponding data set and an encoding of status information for the corresponding data set;

(C) sequentially reading the entries of the first set;

(D) transmitting requests via the network to the computers for those identified ones of the data sets that have corresponding entries of the first set that meet predefined status-based selection criteria; and

(E) creating, in response to receiving the identified ones of the data sets, new entries in said first and second sets corresponding to each of at least a subset of the addresses in the received data sets for which there is no corresponding entry in the second set.

- 29 -

1 17. The method of claim 16, wherein said step B includes
2 storing in each of the entries in the second set an address
3 of a corresponding entry in the first set, said second set
4 of entries being for indexing the first set of entries.

1 18. The method of claim 16, including performing steps C,
2 D and E in multiple threads during overlapping time
3 periods, such that while some of the threads are fetching
4 data sets, other ones of the data sets are analyzing
5 fetched data sets.

1 19. The method of claim 18, including:
2 defining a mutex; and
3 each of said threads requesting and waiting for
4 the mutex before accessing the first and second sets
5 of entries in the first and second memories.

1 20. The method of claim 19, including:
2 defining an input buffer and an append buffer in
3 said second memory;
4 storing blocks of sequentially ordered entries
5 from the first set of entries into the input buffer;
6 said sequentially reading step comprising the
7 step of sequentially reading the entries in the input
8 buffer and determining which of said input buffer
9 entries meet said predefined status-based selection
10 criteria;
11 storing in said append buffer all entries to be
12 added to said first memory; and
13 moving multiple entries in the append buffer to
14 the first memory.

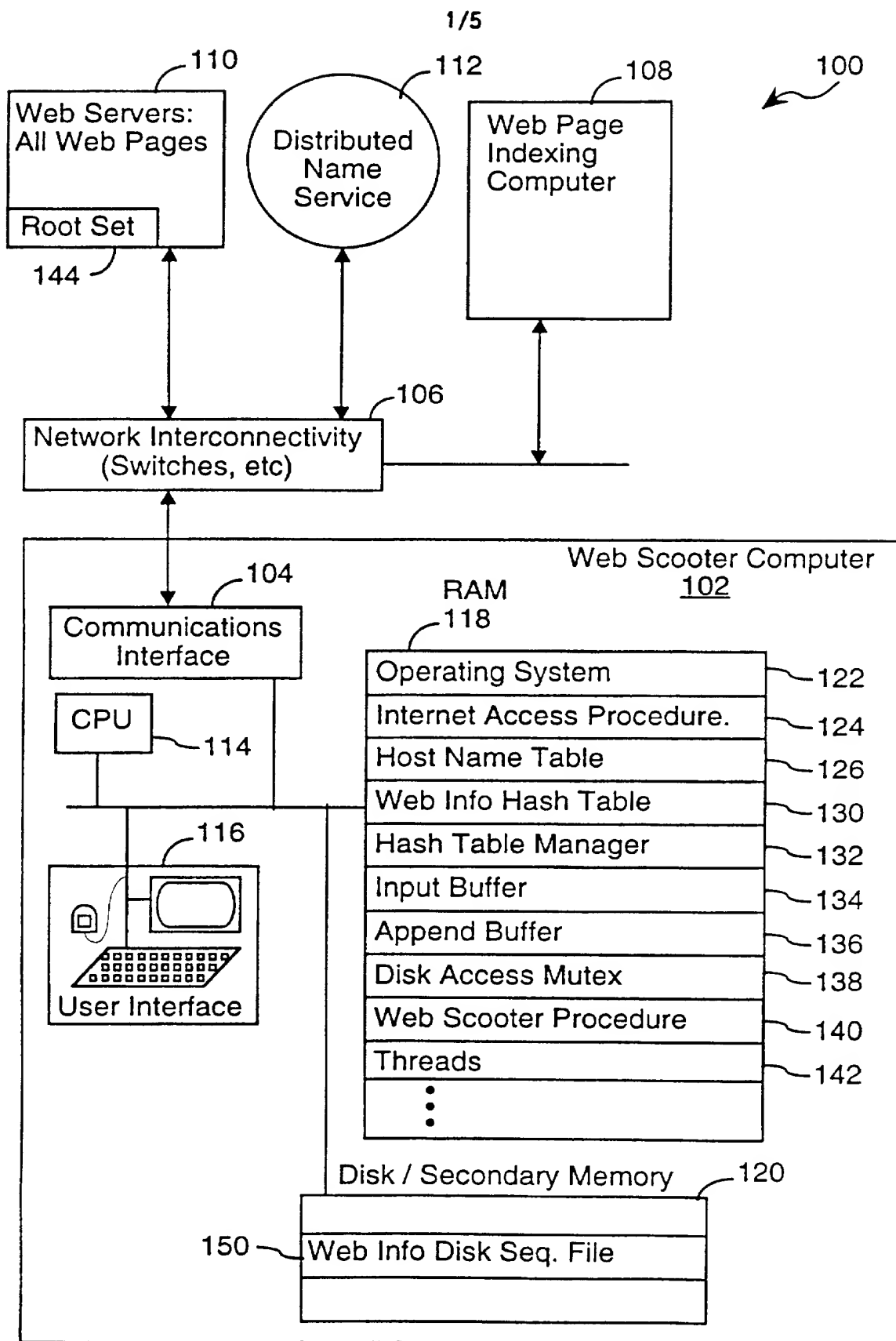


FIG. 1

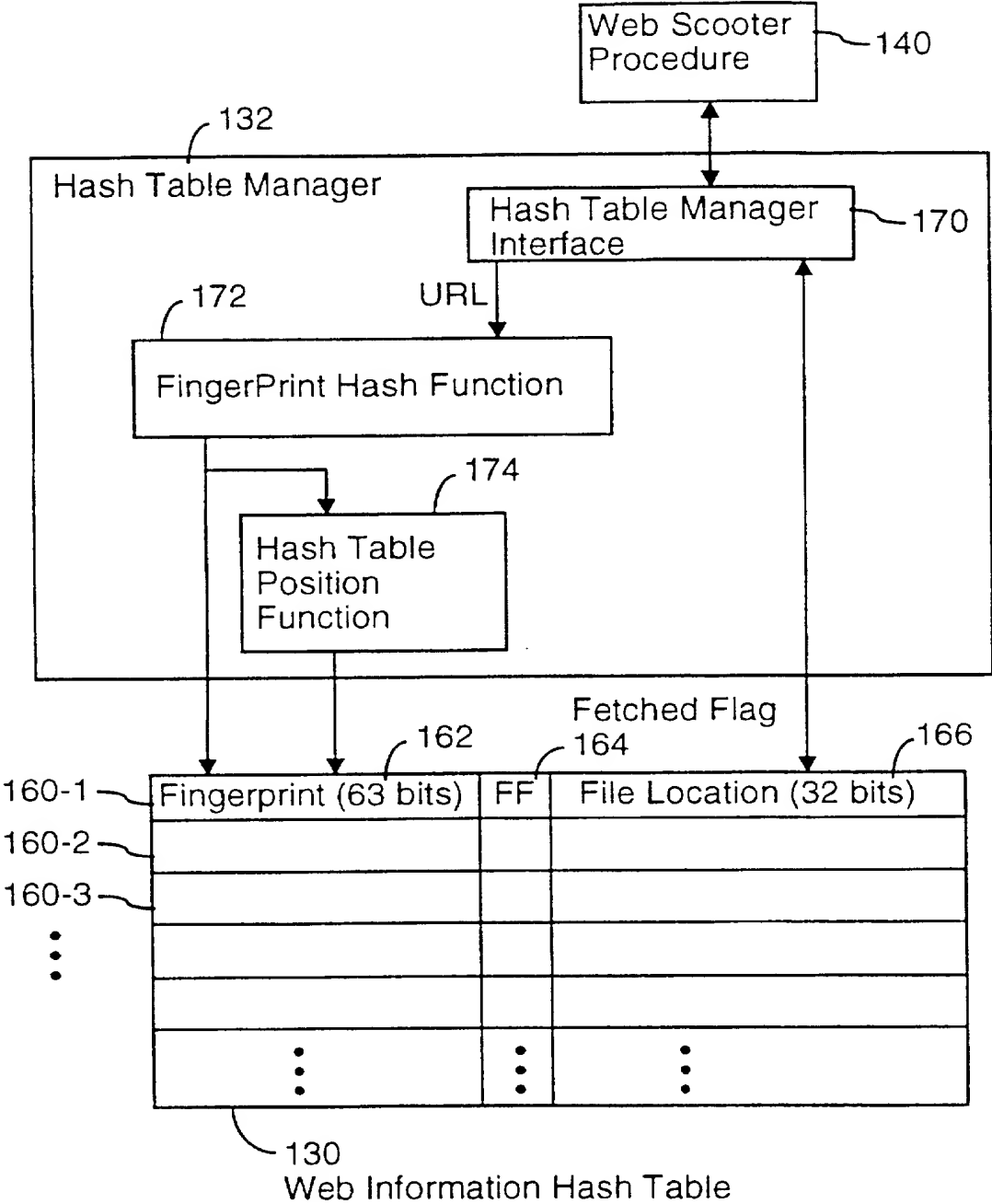


FIG. 2

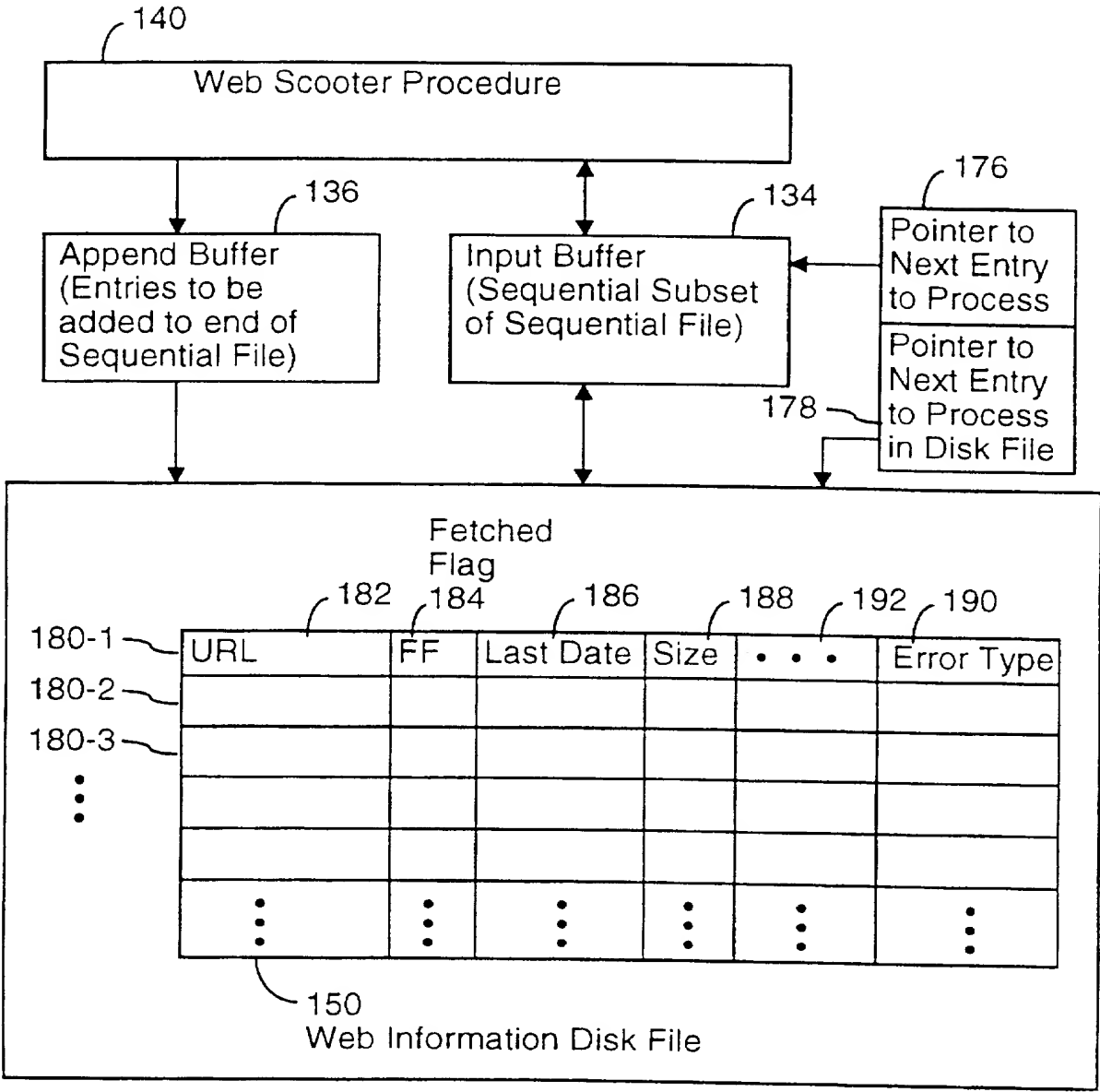


FIG. 3

4/5

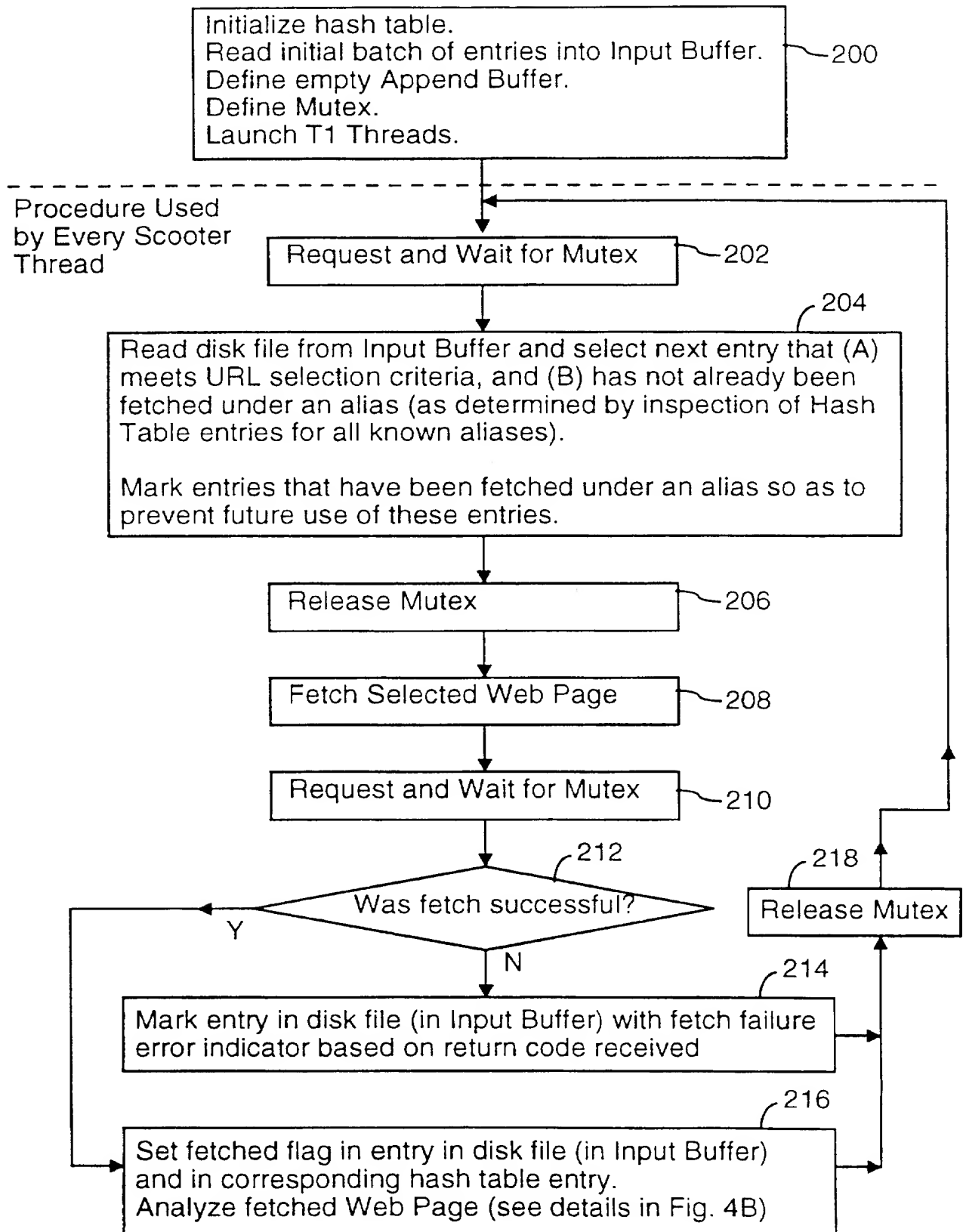


FIG. 4A

5/5

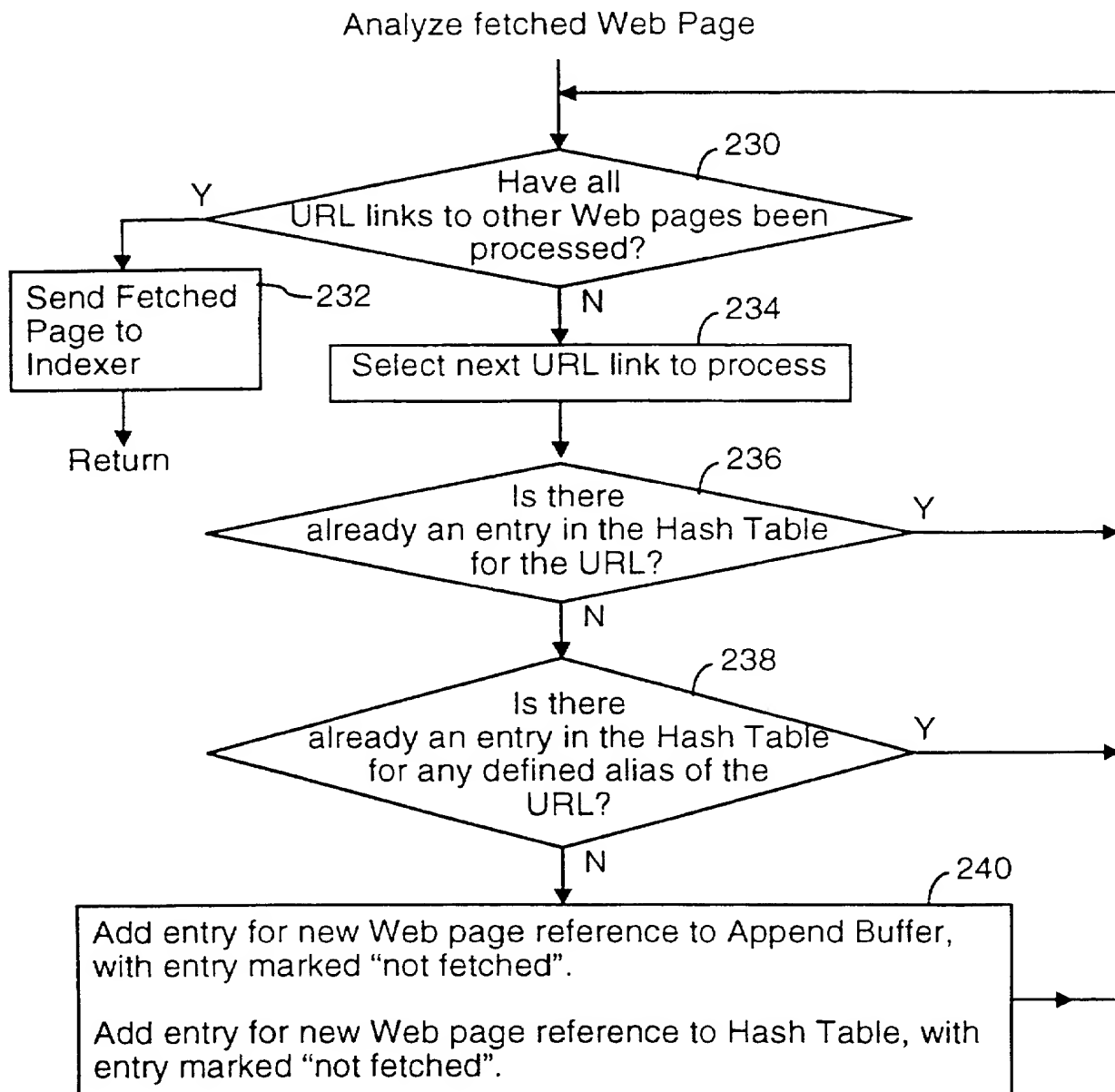


FIG. 4B

FIG. 4A

FIG. 4B

FIG. 4

INTERNATIONAL SEARCH REPORT

International Application No
PCT/US 96/19831

A. CLASSIFICATION OF SUBJECT MATTER
IPC 6 G06F17/30

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)
IPC 6 G06F

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	PERSONAL COMPUTER MAGAZINE, JAN. 1996, VNU BUSINESS PUBLICATIONS, UK, pages 90-92, 94, 97 - 98, 100, XP000646762 SIMPSON D ET AL: "The searchers World Wide Web search engines" see the whole document ---	1,6,11, 16
A	ELECTRONIC LIBRARY, OCT. 1995, LEARNED INFORMATION, UK, vol. 13, no. 5, ISSN 0264-0473, pages 467-476, XP000647883 SHA V T: "Cataloguing Internet resources: the library approach" see the whole document -----	1,6,11, 16

☐ Further documents are listed in the continuation of box C.

☐ Patent family members are listed in annex.

* Special categories of cited documents:

- "A" document defining the general state of the art which is not considered to be of particular relevance
- "E" earlier document but published on or after the international filing date
- "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- "O" document referring to an oral disclosure, use, exhibition or other means
- "P" document published prior to the international filing date but later than the priority date claimed

- "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
- "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
- "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
- "&" document member of the same patent family

Date of the actual completion of the international search

4 April 1997

Date of mailing of the international search report

16.04.97

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Authorized officer

Katerbau, R